

University of Dundee

DataMirror: Reflecting on One's Data Self

Htait, Amal; Azzopardi, Leif; Nicol, Emma; Moncur, Wendy

Published in:

Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)

DOI:

[10.1145/3397271.3401398](https://doi.org/10.1145/3397271.3401398)

Publication date:

2020

Document Version

Peer reviewed version

[Link to publication in Discovery Research Portal](#)

Citation for published version (APA):

Htait, A., Azzopardi, L., Nicol, E., & Moncur, W. (2020). DataMirror: Reflecting on One's Data Self: (A Tool for Social Media Users to explore their Digital Footprints). In Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20) (pp. 2125-2128). (SIGIR 2020 - Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval). Association for Computing Machinery (ACM). <https://doi.org/10.1145/3397271.3401398>

General rights

Copyright and moral rights for the publications made accessible in Discovery Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from Discovery Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain.
- You may freely distribute the URL identifying the publication in the public portal.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

DataMirror: Reflecting on One's Data Self

A Tool for Social Media Users to Explore Their Digital Footprint

Amal Htait
amal.htait@strath.ac.uk
University of Strathclyde
Glasgow, UK

Emma Nicol
enicol001@dundee.ac.uk
University of Dundee
Dundee, UK

Leif Azzopardi
leif.azzopardi@strath.ac.uk
University of Strathclyde
Glasgow, UK

Wendy Moncur
wmoncur@dundee.ac.uk
University of Dundee
Dundee, UK

ABSTRACT

Small pieces of data that are shared online, over time and across multiple social networks, have the potential to reveal more cumulatively than a person intends. This could result in harm, loss or detriment to them depending what information is revealed, who can access it, and how it is processed. But how aware are social network users of how much information they are actually disclosing? And if they could examine all their data, what cumulative revelations might be found that could potentially increase their risk of various online threats (social engineering, fraud, identity theft, loss of face, etc.)? In this paper, we present DataMirror, an initial prototype tool, that enables social network users to aggregate their online data so that they can search, browse and visualise what they have put online. The aim of the tool is to investigate and explore people's awareness of their data self that is projected online; not only in terms of the volume of information that they might share, but what it may mean when combined together, what pieces of sensitive information may be gleaned from their data, and what machine learning may infer about them given their data.

KEYWORDS

Information Revelation; Digital Identity; Data Self; Privacy; Security

ACM Reference Format:

Amal Htait, Leif Azzopardi, Emma Nicol, and Wendy Moncur. 2020. DataMirror: Reflecting on One's Data Self: A Tool for Social Media Users to Explore Their Digital Footprint. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*, July 25–30, 2020, Virtual Event, China. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3397271.3401398>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '20, July 25–30, 2020, Virtual Event, China

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-8016-4/20/07...\$15.00

<https://doi.org/10.1145/3397271.3401398>

1 INTRODUCTION

Online Social Networks (OSNs) allow people to build a network of connections with others in addition to the ability to construct personalised profiles and post content. Such profiles and posts present an opportunity for people to introduce themselves while expressing their personality, thoughts, feelings, and other personal data on the web (e.g., interests, opinions, livelihood, place of work, relationship status, sexual orientation, religion, etc. [5, 9, 10]).

The small pieces of information shared online across multiple networks, may seem innocuous or harmless, individually. However, over time they may reveal more cumulatively than the person intends, as these identifiable traces can be linked together and exploited. For example, many posts may reveal that the person lives alone, while their jogging data shows the routes that they take and where they live, while other posts taken together might suggest, via machine learning tools, that the person is depressed. Thus, these shared data can lead to revealing more about one's identity, habits, work/life patterns, personality, and so forth than they intend – which may result in privacy exposure risks. Such risks can have potentially negative and even disastrous consequences for the person (e.g., identity theft [1], financial losses, imprisonment, damage of reputation [2]), for their employer (e.g., by creating opportunities for cyber-crime, damage to corporate reputation, etc.), and even for national security [3].

How can people use social media, and enjoy its benefits, while minimising their risk of negative or unintended consequences? One possible solution proposed in [3, 7] is the use of a personal informatics system, that enables people to examine and reflect upon the details that they are sharing online, so as to increase their awareness of the privacy risks. For example, the *DataSelfie*¹ project provided a browser plug for Facebook to show individuals what their interactions online might reveal about them. Another project called *WASP* [8] provided a prototype of a personal web archive and search system, integrating archiving, indexing, and reproduction technology into a single application. The prototype offered a customised search engine which enabled users to explore the pages they had visited online.

¹<http://dataselfie.it/>

In this work we are concerned with mapping, presenting and visualising people's profiles and posts so that they can examine and explore information they have explicitly shared online (with other people). Given the recent legislation by the European Union, e.g. General Data Protection Regulation (GDPR)², end-users can request a copy of their own data. However, such bulk downloads are in machine readable format which is not easy to explore. Thus, our tool aims to provide a mechanism for people to explore their downloaded data across multiple platforms. The presented tool, DataMirror, is an initial prototype that will be used as a starting point to study how people reflect upon their own data sharing practices and behaviours, as well as learning more about how their data may be used, combined, and potentially exploited, in different ways – which can lead to security and privacy risks.

DataMirror is a self-contained application (which can be installed via Docker). It is designed to run on the person's machine as opposed to a web service, because the data to be used is their own personal data, which will invariably contain sensitive and personal information about the person.

aims to provide a way to start engaging and motivating people to think about different questions that they might have about their data, and to start to highlight what, given their data, someone or some system, might be able to learn about them, and how that might put them at risk or expose them in some way. To provide different views on their data, users can search and browse through their data, as well as interact with the visualisations.

Identifying Personal Data. Housed on profile pages and disclosed through posts, people reveal different pieces of content that help to personally identify the person (which could be used to compromise their identity [1]). Figure 2 shows an example of extracted details across the different networks.

Figure 2: User’s personal data, collected from different OSNs, presented in a same tree diagram to help the user visualise his cumulated information.

Figure 3: A word-cloud of the most frequent words in a user's online posts, clustered by Topic, reflecting a mainly *political* topic and a use of a very simple vocabulary.

³JavaScript Object Notation

⁴A search engine based on the Lucene library: <https://www.elastic.co/>

⁵<https://www.elastic.co/what-is/kibana>

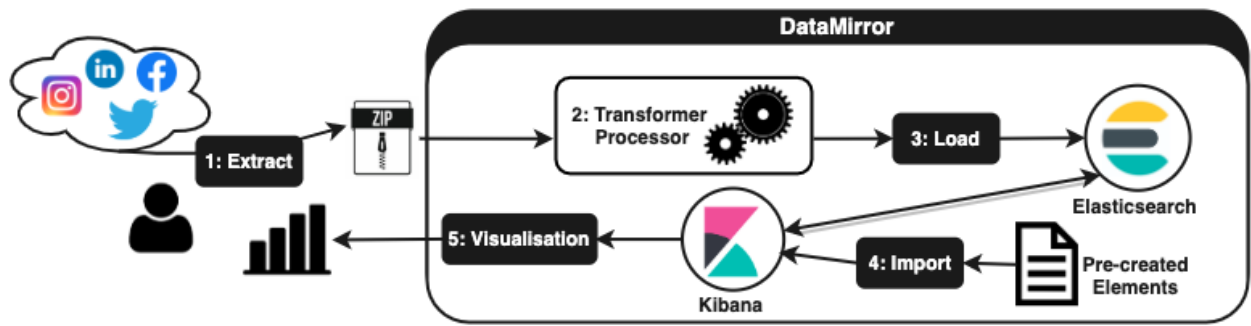


Figure 1: The DataMirror Tool Architecture where the user download their data from the various social media sites for DataMirror to process it in order to provide different insights and visualisations.

Activity. The volume and timing of users' online activities may also provide different insights into their behaviour. Figure 4 aims to show users their online activities (e.g., comments, posts) on different OSNs, divided into time intervals, so that they can observe the potential patterns – which others could potentially glean. For example, employers may be able to observe that staff are spending way too much time on social media during work hours.

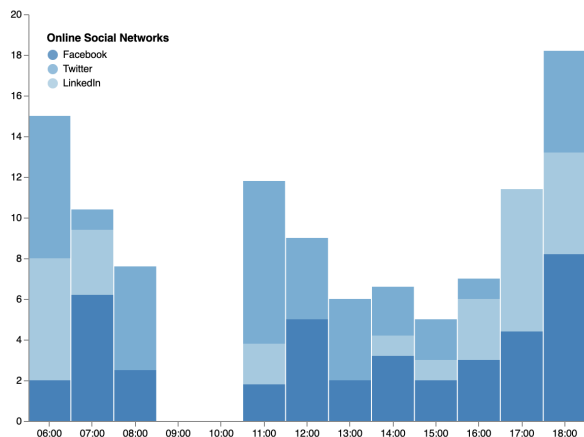


Figure 4: An example of a user's online activity on different OSNs presented in bar diagram, divided into time intervals.

Location. The shared geolocation data in online posts (e.g., check-in) can also reveal more than intended about the user, which could be used to infer their home address, place of work, places they frequent i.e. gyms, shops, parks, cafes. Figure 5 aims to show the users what location information they are sharing which could be used by those with malign intent.

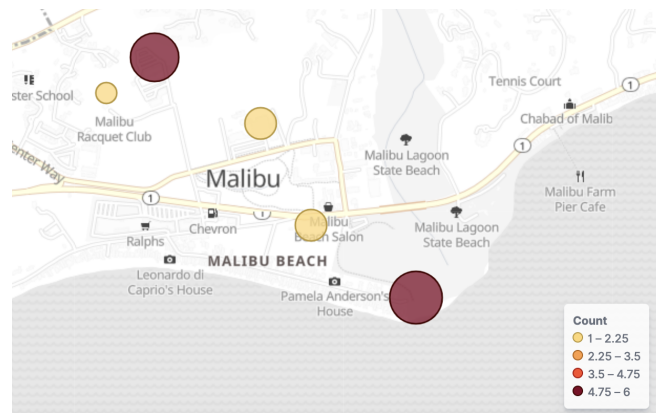


Figure 5: An example of a geolocation map representing the locations posted by the user, with larger/darker dots for the more frequently visited locations.

Connections. The connections that users create can also provide other vulnerabilities e.g. by providing other paths to reveal information about them, by connecting them to minority groups (that could trigger cyber-bullying acts against them), and who they communicate with the most/least in their connections. Figure 7 shows a name-based word cloud where the larger the name the more connections between them.



Figure 7: Word-cloud representing the user's connections: the size of name increases with the closeness of the friend.

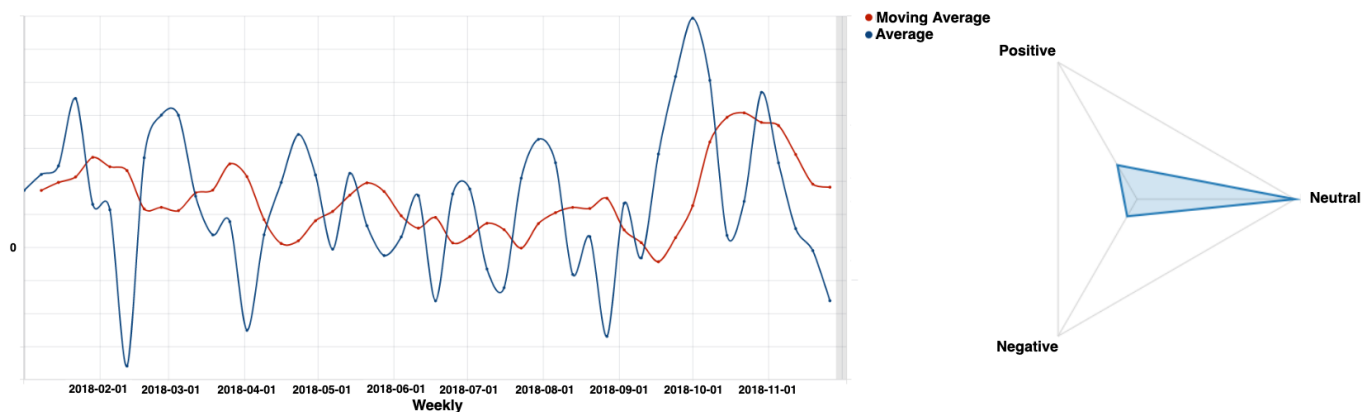


Figure 6: An example of Sentiment Detection Visualisations within the online posts of an anonymous user, with an average of neutral sentiment, and a tendency for subjective posts to be more positive than negative.

Sentiment. Figure 6 shows how posts over time, and collectively, can be interpreted in terms of sentiment, providing potential cues regarding a user’s state of mind and outlook. Sentiment, taken together with interests, could provide information that potential employers might use when making hiring decisions [4], or by others with malicious intent, to socially engineer situations to gain more information or to infiltrate the user’s data further.

3 SUMMARY AND FUTURE WORK

Our initial prototype provides a tool for searching through and visualising one’s data from various social media sites. It lets people explore their own data, it shows them what high level summaries generated from their data provides and what their overall online presence projects to others. Over the course of the project, we aim to extend the prototype to include additional visualisations and deeper content analysis, for example, to identify content as being sensitive (phone numbers, addresses, account numbers, passwords, etc.), to classify content as being personal to the user, as well as including other machine learning algorithms that can provide ratings of people’s mood, level of depression, personality, etc., in addition to including multi-media content (e.g., pictures, videos). Our research will then investigate, through user-studies, what people want to find out from their own data, what questions they want to ask of it, and how we can best present and visualise that to them, in order to promote better awareness and understanding of the risks and potential consequences that sharing many small pieces of information can have cumulatively – as part of the Cumulative Revelations of Personal Data project. The prototype is currently available via GitHub, see <https://github.com/cumulative-revelations/DataMirror>.

ACKNOWLEDGMENTS

Cumulative Revelations of Personal Data This project is supported by the UKRI’s EPSRC under Grant Numbers: EP/R033889/1 and EP/R033854/1.

REFERENCES

- [1] Alessandro Acquisti and Ralph Gross. 2009. Predicting social security numbers from public data. *Proceedings of the National academy of sciences* 106, 27 (2009), 10975–10980.
- [2] Hongliang Chen, Christopher E. Beaudoin, and Traci Hong. 2016. Protecting oneself online: The effects of negative privacy experiences on privacy protective behaviors. *Journalism and Mass Communication Quarterly* 93, 2 (2016), 409–429.
- [3] Judson C Dressler, Christopher Bronk, and Daniel S Wallach. 2015. Exploiting military OpSec through open-source vulnerabilities. In *MILCOM 2015-2015 IEEE Military Communications Conference*. IEEE, 450–458.
- [4] Evanthis Faliagka, Lazaros Iliadis, Ioannis Karydis, Maria Rigou, Spyros Sioutas, Athanasios Tsakalidis, and Giannis Tzimas. 2014. On-line consistent ranking on e-recruitment: seeking the truth behind a well-formed CV. *Artificial Intelligence Review* 42, 3 (2014), 515–528.
- [5] Oliver L. Haimson, Jed R. Brubaker, Lynn Dombrowski, and Gillian R. Hayes. 2016. Digital footprints and changing networks during online identity transitions. In *Conference on Human Factors in Computing Systems - Proceedings*. Association for Computing Machinery, 2895–2907.
- [6] Clayton J Hutto and Eric Gilbert. 2014. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Eighth international AAAI conference on weblogs and social media*.
- [7] Danesh Irani, Steve Webb, Kang Li, and Calton Pu. 2011. Modeling unintended personal-information leakage from multiple online social networks. *IEEE Internet Computing* 15, 3 (2011), 13–19.
- [8] Johannes Kiesel, Arjen P de Vries, Matthias Hagen, Benno Stein, and Martin Potthast. 2018. WASP: web archiving and search personalized. (2018).
- [9] Hanna Krasnova, Sarah Spiekermann, Ksenia Koroleva, and Thomas Hildebrand. 2010. Online social networks: Why we disclose. *Journal of Information Technology* 25, 2 (jun 2010), 109–125.
- [10] Ning Xia, Han Hee Song, Yong Liao, Marios Iliofotou, Antonio Nucci, Zhi-Li Zhang, and Aleksandar Kuzmanovic. 2013. *Mosaic: Quantifying Privacy Leakage in Mobile Networks*. 564 pages.